

On the complexity of piecewise affine system identification

Fabien Lauer

Université de Lorraine, CNRS, LORIA, UMR 7503, F-54506 Vandœuvre-lès-Nancy, France

September 9, 2015

Abstract

The paper provides results regarding the computational complexity of hybrid system identification. More precisely, we focus on the estimation of piecewise affine (PWA) maps from input-output data and analyze the complexity of computing a global minimizer of the error. Previous work showed that a global solution could be obtained for continuous PWA maps with a worst-case complexity exponential in the number of data. In this paper, we show how global optimality can be reached for a slightly more general class of possibly discontinuous PWA maps with a complexity only polynomial in the number of data, however with an exponential complexity with respect to the data dimension. This result is obtained via an analysis of the intrinsic classification subproblem of associating the data points to the different modes. In addition, we prove that the problem is NP-hard, and thus that the exponential complexity in the dimension is a natural expectation for any exact algorithm.

1 Introduction

Hybrid system identification aims at estimating a model of a system switching between different operating modes from input-output data. More precisely, most of the literature considers autoregressive with external input (ARX) models to cast the problem as a regression one [1]. Then, two cases can be distinguished: switching regression, where the system arbitrarily switches from one mode to another, and piecewise affine (PWA) regression, where the switches depend on the regressors. A number of methods with satisfactory performance in practice are now available for these problems [2]. However, compared with linear system identification, a major weakness of these methods is their lack of guarantees.

For the particular case of noiseless data, the algebraic method [3] provides a solution to switching regression with a small number of modes. However, the quality of the estimates quickly degrades with the increase of the noise level. A few sparsity-based methods [4, 5] also offer guarantees in the noiseless case, but these are subject to a condition on both the data and the sought solution. In the presence of noise, most methods consider the minimization of the error of the model over the data [1]. While this does not necessarily yields the best predictive model (due to issues like identifiability, persistence of excitation and access to a limited amount of data), obtaining statistical guarantees with such an approach has a long history in statistics and system identification [6]. However, such results are not available for hybrid systems. This is probably due to the fact that minimizing the error of a hybrid model is a difficult nonconvex optimization problem involving the simultaneous classification of the data points into modes and the regression of a submodel for each mode. Thus, theoretical guarantees could only be obtained under the rather strong assumption that this problem has been solved to global optimality and most of the literature [7, 8, 9, 10, 11, 12] focuses on this issue with heuristics of various degrees of accuracy and computational efficiency. Many recent works [4, 13, 14, 15, 16, 5] try to avoid local minima by considering convex formulations, but these only yield optimality with respect to a relaxation of the original problem. Global optimality in the presence of noise was only reached in [17] for a particular class of continuous PWA maps known as hinging-hyperplanes by reformulating the problem as a mixed-integer program solved by

branch-and-bound techniques. However, such optimization problems are NP-hard [18] and branch-and-bound algorithms have a worst-case complexity exponential in the number of integer variables, here proportional to the number of data and the number of modes.

Inspired by related clustering problems, such as the minimization of the sum of squared distances between points and their group centers, we could minimize the hybrid model error by enumerating all possible classifications of the points. But the number of classifications is exponential in the number of data. Conversely, the other approach enumerating a sample of values for the real variables of the problem is exponential in the dimension and can only offer an approximate solution.

Overall, the literature does not provide a method that can guarantee both the optimality and the computability of a global minimizer of the error, while the computational complexity of this problem remains unknown and cannot be deduced from the NP-hardness of classical clustering problems [19] (see [18, 20] for an introduction to computational complexity and its relevance to control theory).

Contribution The paper provides two results regarding the computational complexity of PWA regression, and more precisely for the problem of finding a global minimizer of the error of a PWA model, formalized in Sect. 2. First, we show in Sect. 3 that the problem is NP-hard. Then, we show in Sect. 4 that, for any fixed dimension of the data, an exact solution can be computed in time polynomial in the number of data via an enumeration of all possible classifications. To obtain this result and avoid the exponential growth of the number of classifications with the number of data, we show that, in PWA regression, the classification of the data points is highly constrained and the number of classifications to test can be limited. The price to pay for this gain is an exponential complexity with respect to the data dimension and the number of modes. Future work is outlined in Sect. 5.

Notations We use the indicator function $\mathbf{1}_E$ of an event E that is 1 if the event occurs and 0 otherwise. We define $\text{sign}(u) = 1$ if $u \geq 0$ and -1 otherwise. Given a set of labels $\mathcal{Q} \subset \mathbb{Z}$ and a set of N points, a labeling of these points is any $\mathbf{q} \in \mathcal{Q}^N$. We use $j = \arg\max_{k \in \mathcal{Q}} u(k)$ as a shorthand for $j = \min\{l \in \arg\max_{k \in \mathcal{Q}} u(k)\}$. Given two sets, \mathcal{X} and \mathcal{Y} , $\mathcal{Y}^{\mathcal{X}}$ is the set of functions from \mathcal{X} to \mathcal{Y} .

2 Problem formulation

As in most works, we concentrate on discrete-time PWARX system identification considered as a PWA regression problem with regression vectors $\mathbf{x}_i = [y_{i-1}, \dots, y_{i-n_y}, u_i, \dots, u_{i-n_u}]^T \in \mathcal{X}$ built from past inputs u_{i-k} and outputs y_{i-k} . Since we are interested in computational complexity results, we restrain the data to rational, digitally representable, values and set $\mathcal{X} \subseteq \mathbb{Q}^d$. The outputs are assumed to be generated by a PWA system f as $y_i = f(\mathbf{x}_i) + v_i$, where v_i is a noise term. More precisely, PWA models can be expressed via a set of n affine submodels and a function $h : \mathcal{X} \rightarrow \mathcal{Q} = \{1, \dots, n\}$ determining the active submodel: $f(\mathbf{x}) = \mathbf{w}_{h(\mathbf{x})}^T \bar{\mathbf{x}}$, where $\bar{\mathbf{x}} = [\mathbf{x}^T, 1]^T$.

We call the function h a classifier as it classifies the data points in the different modes. Typically, PWA systems are defined with h implementing a polyhedral partition of \mathcal{X} , with modes possibly spanning unions of polyhedra. However, in most of the literature on PWA system identification [1, 7, 8, 9, 16], h is estimated within the family of linear classifiers

$$\mathcal{H} = \{h \in \mathcal{Q}^{\mathcal{X}} : h(\mathbf{x}) = \arg\max_{k \in \mathcal{Q}} \mathbf{h}_k^T \mathbf{x} + b_k, \mathbf{h}_k \in \mathbb{Q}^d, b_k \in \mathbb{Q}\}, \quad (1)$$

based on a set of n linear functions and for which a mode spanning a union of polyhedra must be modeled as several modes with similar affine submodels. For PWA maps with $n = 2$ modes, h is a binary classifier for which it is common to consider its output in $\mathcal{Q} = \{-1, +1\}$ instead of $\{1, 2\}$. Such a binary classifier can be obtained by taking the sign of a real-valued function. If this function is linear (or affine), then we obtain a linear classifier, which is equivalent to a separating hyperplane dividing the input space \mathcal{X} in two half-spaces. In this case, the function class \mathcal{H} can be defined as

$$\mathcal{H} = \{h \in \mathcal{Q}^{\mathcal{X}} : h(\mathbf{x}) = \text{sign}(\mathbf{h}^T \mathbf{x} + b), \mathbf{h} \in \mathbb{Q}^d, b \in \mathbb{Q}\} \quad (2)$$

with a single set of parameters (\mathbf{h}, b) corresponding to the normal to the hyperplane and the offset from the origin. An equivalence with the multi-class formulation in (1) is obtained by using $\mathbf{h} = \mathbf{h}_1 - \mathbf{h}_2$ and $b = b_1 - b_2$.

In this paper, we consider the common estimation approach of minimizing the error on N data pairs $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathbb{Q}$, measured pointwise by a loss function $\ell : \mathbb{Q} \rightarrow \mathbb{Q}^+$ as

$$\ell(y_i - f(\mathbf{x}_i)) = \sum_{j \in \mathcal{Q}} \mathbf{1}_{h(\mathbf{x}_i)=j} \ell(y_i - \mathbf{w}_j^T \bar{\mathbf{x}}_i).$$

More precisely, we focus on well-posed instances of the problem where N is significantly larger than the dimension d and the number of modes n is given. Indeed, with free n the problem is ill-posed as the solution is only defined up to a trade-off between the number of modes and the model accuracy. For the converse well-posed approach that minimizes n for a given error bound, a complexity analysis can be found in [21]. Under these assumptions, the problem is as follows.

Problem 1 (Error-minimizing PWA regression). *Given a data set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N \in (\mathcal{X} \times \mathbb{Q})^N$ with $\mathcal{X} \subseteq \mathbb{Q}^d$ and an integer $n \in [2, N/(d+1)]$, find a global solution to*

$$\min_{\mathbf{w} \in \mathbb{Q}^{n(d+1)}, h \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \sum_{j \in \mathcal{Q}} \mathbf{1}_{h(\mathbf{x}_i)=j} \ell(y_i - \mathbf{w}_j^T \bar{\mathbf{x}}_i), \quad (3)$$

where $\mathbf{w} = (\mathbf{w}_j)_{j \in \mathcal{Q}}$ is the concatenation of all parameter vectors and $\mathcal{H} \subset \mathcal{Q}^{\mathcal{X}}$ is the set of n -category linear classifiers as in (1) or (2).

The following analyzes the time complexity of Problem 1 under the classical model of computation known as a Turing machine [18]. The time complexity of a problem is the lowest time complexity of an algorithm solving any instance of that problem, where the time complexity of an algorithm is the maximal number of steps occurring in the computation of the corresponding Turing machine program. The loss function ℓ is assumed to be computable in polynomial time throughout the paper.

3 NP-hardness

This section contains the proof of the following NP-hardness result, where an NP-hard problem is one that is at least as hard as any problem from the class NP of nondeterministic polynomial time decision problems [18] (NP is the class of all decision problems for which a solution can be certified in polynomial time).

Theorem 1. *With a loss function ℓ such that $\ell(e) = 0 \Leftrightarrow e = 0$, Problem 1 is NP-hard.*

The proof uses a reduction from the partition problem, known to be NP-complete [18], i.e., a problem that is both NP-hard and in NP.

Problem 2 (Partition). *Given a multiset (a set with possibly multiple instances of its elements) of d positive integers, $S = \{s_1, \dots, s_d\}$, decide whether there is a multisubset $S_1 \subset S$ such that*

$$\sum_{s_i \in S_1} s_i = \sum_{s_i \in S \setminus S_1} s_i.$$

More precisely, we will reduce Problem 2 to the decision form of Problem 1.

Problem 3 (Decision form of PWA regression). *Given a data set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N \in (\mathcal{X} \times \mathbb{Q})^N$, an integer $n \in [2, N/(d+1)]$ and a threshold $\epsilon \geq 0$, decide whether there is a pair $(\mathbf{w}, h) \in \mathbb{Q}^{n(d+1)} \times \mathcal{H}$ such that*

$$\frac{1}{N} \sum_{i=1}^N \sum_{j \in \mathcal{Q}} \mathbf{1}_{h(\mathbf{x}_i)=j} \ell(y_i - \mathbf{w}_j^T \bar{\mathbf{x}}_i) \leq \epsilon, \quad (4)$$

where \mathcal{H} is the set of linear classifiers as in (1) or (2) and the loss function ℓ is such that $\ell(e) = 0 \Leftrightarrow e = 0$.

Proposition 1. *Problem 3 is NP-complete.*

Proof. Since given a candidate pair (\mathbf{w}, h) the condition (4) can be verified in polynomial time, Problem 3 is in NP. Then, the proof of its NP-completeness proceeds by showing that the Partition Problem 2 has an affirmative answer if and only if Problem 3 with $\epsilon = 0$ has an affirmative answer.

Given an instance of Problem 2, let $N = 2d + 3$, $n = 2$, $\mathcal{Q} = \{-1, 1\}$ and build a data set with

$$(\mathbf{x}_i, y_i) = \begin{cases} (s_i \mathbf{e}_i, s_i), & \text{if } 1 \leq i \leq d \\ (-s_{i-d} \mathbf{e}_{i-d}, s_{i-d}), & \text{if } d < i \leq 2d \\ (\mathbf{s}, 0), & \text{if } i = 2d + 1 \\ (-\mathbf{s}, 0), & \text{if } i = 2d + 2 \\ (\mathbf{0}, 0), & \text{if } i = 2d + 3, \end{cases}$$

where \mathbf{e}_k is the k th unit vector of the canonical basis for \mathbb{Q}^d and $\mathbf{s} = \sum_{k=1}^d s_k \mathbf{e}_k$. If Problem 2 has an affirmative answer, then, using the notations of (2), we can set

$$\mathbf{w}_1 = \sum_{k \in I_1} \bar{\mathbf{e}}_k - \sum_{k \in I_{-1}} \bar{\mathbf{e}}_k, \quad \mathbf{w}_{-1} = -\mathbf{w}_1, \quad \mathbf{h} = \sum_{k \in I_1} \mathbf{e}_k - \sum_{k \in I_{-1}} \mathbf{e}_k, \quad b = 0,$$

where $\bar{\mathbf{e}}_k = [\mathbf{e}_k^T, 0]^T$, I_1 is the set of indexes of the elements of S in S_1 and $I_{-1} = \{1, \dots, d\} \setminus I_1$. This gives

$$\mathbf{w}_1^T \bar{\mathbf{x}}_i = \begin{cases} s_i = y_i, & \text{if } i \leq d \text{ and } i \in I_1 \\ -s_i, & \text{if } i \leq d \text{ and } i \in I_{-1} \\ s_{i-d} = y_i, & \text{if } i > d \text{ and } i-d \in I_{-1} \\ -s_{i-d}, & \text{if } i > d \text{ and } i-d \in I_1 \\ \sum_{k \in I_1} s_k - \sum_{k \in I_{-1}} s_k = 0 = y_i, & \text{if } i = 2d + 1 \\ \sum_{k \in I_{-1}} s_k - \sum_{k \in I_1} s_k = 0 = y_i, & \text{if } i = 2d + 2 \\ 0 = y_i, & \text{if } i = 2d + 3 \end{cases}$$

and we can similarly show that

$$\mathbf{w}_{-1}^T \bar{\mathbf{x}}_i = y_i, \quad \text{if } i \in I_{-1} \vee i-d \in I_1 \vee i > 2d,$$

while $\mathbf{h}^T \mathbf{x}_i$ is positive if $i \in I_1 \vee i-d \in I_{-1}$ and negative if $i \in I_{-1} \vee i-d \in I_1$. Therefore, for all points, $\mathbf{w}_{h(\mathbf{x}_i)}^T \bar{\mathbf{x}}_i = y_i$, $i = 1, \dots, 2d + 3$, and the cost function of Problem 1 is zero, yielding an affirmative answer for Problem 3.

It remains to prove that if (4) holds with $\epsilon = 0$, then Problem 2 has an affirmative answer. To see this, note that due to ℓ being positive, a zero cost implies a zero loss for all data points. Thus, by $\ell(e) = 0 \Leftrightarrow e = 0$, if (4) holds with $\epsilon = 0$,

$$\mathbf{w}_{h(\mathbf{x}_i)}^T \bar{\mathbf{x}}_i = y_i, \quad i = 1, \dots, 2d + 3. \quad (5)$$

Also note that if $h(\mathbf{x}_i) = h(\mathbf{x}_{i+d}) = 1$ for some $i \leq d$, we have $s_i w_{1,i} + w_{1,d+1} = -s_i w_{1,i} + w_{1,d+1} = s_i$. This is only possible if $s_i = 0$, which is not the case (otherwise we can simply remove s_i without influencing the partition problem), or if $w_{1,d+1} = s_i$. The latter is impossible if $h(\mathbf{x}_i) = h(\mathbf{x}_{i+d})$ since h is a linear classifier that must return the same category for all points on the line segment between \mathbf{x}_i and \mathbf{x}_{i+d} , which includes the origin $\mathbf{x}_{2d+3} = \mathbf{0}$ and thus would imply by (5) that $\mathbf{w}_1^T \bar{\mathbf{x}}_{2d+3} = w_{1,d+1} = y_{2d+3} = 0$. As a consequence, $h(\mathbf{x}_i) = h(\mathbf{x}_{i+d}) = 1$ cannot hold, and since we can similarly show that $h(\mathbf{x}_i) = h(\mathbf{x}_{i+d}) = -1$ cannot hold, we have $h(\mathbf{x}_i) \neq h(\mathbf{x}_{i+d})$ for all $i \leq d$. Hence, (5) leads to

$$\mathbf{w}_1^T \bar{\mathbf{x}}_i \neq y_i \Rightarrow \mathbf{w}_1^T \bar{\mathbf{x}}_{d+i} = y_{d+i}, \quad i = 1, \dots, d \quad (6)$$

$$\mathbf{w}_{-1}^T \bar{\mathbf{x}}_i \neq y_i \Rightarrow \mathbf{w}_{-1}^T \bar{\mathbf{x}}_{d+i} = y_{d+i}, \quad i = 1, \dots, d. \quad (7)$$

Let $\hat{I}_1 = \{i \in \{1, \dots, d\} : \mathbf{w}_1^T \bar{\mathbf{x}}_i = y_i\}$ and $\hat{I}_{-1} = \{1, \dots, d\} \setminus \hat{I}_1$. Then, if $h(\mathbf{x}_{2d+3}) = +1$, $w_{1,d+1} = 0$ and for all $i \leq d$, $\mathbf{w}_1^T \bar{\mathbf{x}}_i = w_i s_i$. Therefore, for all $i \in \hat{I}_1$, $w_i = 1$, while for all $i \in \hat{I}_{-1}$, (6) gives $\mathbf{w}_1^T \bar{\mathbf{x}}_{d+i} = y_{d+i}$, i.e., $-w_i s_i = s_i$ and $w_i = -1$. This leads to

$$\mathbf{w}_1^T \bar{\mathbf{x}}_{2d+1} = \sum_{i \in \hat{I}_1} s_i - \sum_{i \in \hat{I}_{-1}} s_i = -\mathbf{w}_1^T \bar{\mathbf{x}}_{2d+2}.$$

Thus, if $\mathbf{w}_1^T \bar{\mathbf{x}}_{2d+1} = y_{2d+1} = 0$ or $\mathbf{w}_1^T \bar{\mathbf{x}}_{2d+2} = y_{2d+2} = 0$, a valid partition in the sense of Problem 2 is obtained with $S_1 = \{s_i\}_{i \in \hat{I}_1}$. In addition, if $\mathbf{w}_1^T \bar{\mathbf{x}}_{2d+1} \neq 0$ and $\mathbf{w}_1^T \bar{\mathbf{x}}_{2d+2} \neq 0$, then by (5), $\mathbf{w}_{-1}^T \bar{\mathbf{x}}_{2d+1} = \mathbf{w}_{-1}^T \bar{\mathbf{x}}_{2d+2} = 0$, which by construction implies that $w_{-1,d+1} = 0$. In this case, we redefine $\hat{I}_{-1} = \{i \in \{1, \dots, d\} : \mathbf{w}_{-1}^T \bar{\mathbf{x}}_i = y_i\}$ and $\hat{I}_1 = \{1, \dots, d\} \setminus \hat{I}_{-1}$ to obtain $w_{-1,i} = 1$ for all $i \in \hat{I}_{-1}$ and $w_{-1,i} = -1$ for all $i \in \hat{I}_1$, resulting also in a valid partition by the fact that $\mathbf{w}_{-1}^T \bar{\mathbf{x}}_{2d+2} = \sum_{i \in \hat{I}_1} s_i - \sum_{i \in \hat{I}_{-1}} s_i = 0$. Since a similar reasoning applies to the case $h(\mathbf{x}_{2d+3}) = -1$ by symmetry (substituting \mathbf{w}_{-1} for \mathbf{w}_1), a zero cost, i.e., (4) with $\epsilon = 0$, always implies an affirmative answer to Problem 2. \square

Proof of Theorem 1. Since the decision form of Problem 1 with $\ell(e) = 0 \Leftrightarrow e = 0$, i.e., Problem 3, is NP-complete, Problem 1 with such a loss function is NP-hard (solving Problem 1 also yields the answer to Problem 3 and thus it is at least as hard as Problem 3). \square

4 Polynomial complexity in the number of data

We now state the result regarding the polynomial complexity of Problem 1 with respect to N under the following assumptions, the first of which holds almost surely for randomly drawn data points, while the second one holds for instance for $\ell(e) = e^2$ with a linear time complexity $T(N) = \mathcal{O}(N)$ [22].

Assumption 1. The points $\{\mathbf{x}_i\}_{i=1}^N$ are in general position, i.e., no hyperplane of \mathbb{Q}^d contains more than d points.

Assumption 2. Given $\{(\mathbf{x}_i, y_i)\}_{i=1}^N \in (\mathcal{X} \times \mathbb{Q})^N$, the problem $\min_{\mathbf{v} \in \mathbb{Q}^{d+1}} \sum_{i=1}^N \ell(y_i - \mathbf{v}^T \bar{\mathbf{x}}_i)$ has a polynomial time complexity $T(N)$ for any fixed integer $d \geq 1$.

Theorem 2. For any fixed number of modes n and dimension d , under Assumptions 1–2, the time complexity of Problem 1 is no more than polynomial in the number of data N and in the order of $T(N) \mathcal{O}(N^{dn(n-1)/2})$.

The proof of Theorem 2 relies on the existence of exact algorithms with complexity polynomial in N for the binary case ($n = 2$, Proposition 4) and the multi-class case ($n \geq 3$, Corollary 1). These algorithms are based on a reduction of Problem 1 to a combinatorial search in two steps. The first step reduces the problem to a classification one. Indeed, Problem 1 can be reformulated as the search for the classifier h , since by fixing h , the optimal parameter vectors $\{\mathbf{w}_j\}_{j \in \mathcal{Q}}$ can be obtained by solving n independent linear regression problems on the subsets of data resulting from the classification by h , which, by Assumption 2, can be performed in the polynomial time $T(N)$. This yields the following reformulation of the problem.

Proposition 2. Problem 1 is equivalent to

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{Q}^{n(d+1)}, h \in \mathcal{H}} \quad & \frac{1}{N} \sum_{i=1}^N \sum_{j \in \mathcal{Q}} \mathbf{1}_{h(\mathbf{x}_i)=j} \ell(y_i - \mathbf{w}_j^T \bar{\mathbf{x}}_i) \\ \text{s.t. } \forall j \in \mathcal{Q}, \quad & \mathbf{w}_j \in \operatorname{argmin}_{\mathbf{v} \in \mathbb{Q}^{d+1}} \sum_{i=1}^N \mathbf{1}_{h(\mathbf{x}_i)=j} \ell(y_i - \mathbf{v}^T \bar{\mathbf{x}}_i). \end{aligned} \quad (8)$$

The second step reduces the estimation of h to a combinatorial problem solved in $\mathcal{O}(N^{dn(n-1)/2})$ operations, as detailed in Sect. 4.1–4.2 for $n = 2$ and in Sect. 4.3 for $n \geq 3$.

4.1 Finding the optimal classification

We reduce the complexity of searching for the classifier by considering all possible linear classifications instead of all possible linear classifiers. In other words, we project the class \mathcal{H} of classifiers onto the set of points $S = \{\mathbf{x}_i\}_{i=1}^N$ to reduce a continuous search to a combinatorial problem. This is in line with the techniques used in statistical learning theory [23] for the different purpose of computing error bounds for infinite function classes. Thus, we introduce definitions from this field.

Definition 1 (Projection onto a set). *The projection of a set of classifiers $\mathcal{H} \subset \mathcal{Q}^{\mathcal{X}}$ onto $S = \{\mathbf{x}_i\}_{i=1}^N$, denoted \mathcal{H}_S , is the set of all labelings of S that can be produced by a classifier in \mathcal{H} :*

$$\mathcal{H}_S = \{(h(\mathbf{x}_1), \dots, h(\mathbf{x}_N)) : h \in \mathcal{H}\} \subseteq \mathcal{Q}^N.$$

Definition 2 (Growth function). *The growth function $\Pi_{\mathcal{H}}(N)$ of \mathcal{H} at N is the maximal number of labelings of N points that can be produced by classifiers from \mathcal{H} :*

$$\Pi_{\mathcal{H}}(N) = \sup_{S \in \mathcal{X}^N} |\mathcal{H}_S|.$$

We now focus on binary PWA maps and thus on binary classifiers with output in $\mathcal{Q} = \{-1, +1\}$. For such classifiers, we obviously have $\Pi_{\mathcal{H}}(N) \leq 2^N$ for all N . By further restricting \mathcal{H} to affine classifiers as in (2), results from statistical learning theory (see, e.g., [23]) provide the tighter bound $\Pi_{\mathcal{H}}(N) \leq \left(\frac{eN}{d+1}\right)^{d+1}$, which is polynomial in N and thus promising from the viewpoint of global optimization. However, its proof is not constructive and does not provide an explicit algorithm for enumerating all the labelings. The following theorem, though leading to a looser bound on the growth function, offers a constructive scheme to compute the projection \mathcal{H}_S , which is what we need in order to test all the labelings in \mathcal{H}_S for global optimization.

Theorem 3. *The growth function of the class of binary affine classifiers of \mathbb{Q}^d , \mathcal{H} in (2), is bounded for any $N > d$ by*

$$\Pi_{\mathcal{H}}(N) \leq 2^{d+1} \binom{N}{d} = \mathcal{O}(N^d)$$

and, for any set S of N points in general position, an algorithm builds the projection \mathcal{H}_S in $\mathcal{O}(N^d)$ time.

The proof of Theorem 3 relies on the following proposition, which is illustrated by Fig. 1.

Proposition 3. *For any binary affine classifier h in \mathcal{H} (2) and any finite set of $N > d$ points $S = \{\mathbf{x}_i\}_{i=1}^N$ in general position, there is a subset of points $S_h \subset S$ of cardinality $|S_h| = d$ and a separating hyperplane of parameters $(\mathbf{h}_{S_h}, b_{S_h})$ passing through the points in S_h , i.e.,*

$$\forall \mathbf{x} \in S_h, \quad \mathbf{h}_{S_h}^T \mathbf{x} + b_{S_h} = 0, \quad \text{with } \|\mathbf{h}_{S_h}\| = 1, \quad (9)$$

which yields the same classification of S in the sense that

$$\forall \mathbf{x}_i \in S \setminus S_h, \quad h(\mathbf{x}_i) = \text{sign}(\mathbf{h}_{S_h}^T \mathbf{x}_i + b_{S_h}). \quad (10)$$

Proof sketch. For all classifiers h with separating hyperplanes passing through d points of S , the statement is obvious. For the others passing through p points with $0 \leq p < d$, they can be transformed to pass through additional points without changing the classification of the remaining points. If $p = 0$, it suffices to translate the hyperplane to the closest point. If $0 < p < d$, the hyperplane can be rotated with a plane of rotation that leaves unchanged the subspace spanned by the p points and a minimal angle yielding a rotated hyperplane passing through $p' > p$ points, where $p' \leq d$ by the general position assumption. Iterating this scheme until $p = d$ yields a hyperplane passing through the points in S_h of parameters $(\mathbf{h}_{S_h}, b_{S_h})$ satisfying (9) and (10). \square

We can now prove Theorem 3.

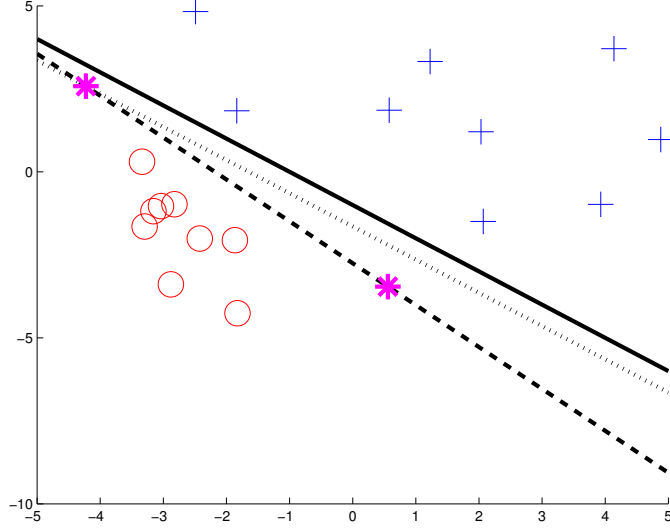


Figure 1: The hyperplane H (plain line) produces the same classification (into $+$ and o) as the hyperplane H_S (dashed line) obtained by a translation (dotted line) and a rotation of H such that it passes through exactly 2 points of S (*).

Proof of Theorem 3. For any labeling \mathbf{q} in \mathcal{H}_S , there is a classifier $h \in \mathcal{H}$ that produces this labeling. Applying Proposition 3 to h , we obtain another classifier h_{S_h} of parameters $(\mathbf{h}_{S_h}, b_{S_h})$ that passes through the points in S_h and that agrees with h on $S \setminus S_h$. Let $\hat{\mathbf{q}} \in \{-1, +1\}^N$ be defined by $\hat{q}_i = h_{S_h}(\mathbf{x}_i)$, $i = 1, \dots, N$. Then, we generate 2^d labelings by setting its entries \hat{q}_i with $i \in S_h$ to all the 2^d combinations of signs (recall that $|S_h| = d$). By construction, there is no labeling of S that agrees with \mathbf{q} on $S \setminus S_h$ other than these 2^d labelings. Since this holds for any $\mathbf{q} \in \mathcal{H}_S$, the cardinality of \mathcal{H}_S cannot be larger than 2^d times the number of hyperplanes passing through d points of S . Since each subset $S_h \subset S$ of cardinality d gives rise to two hyperplanes of opposite orientations, the number of such hyperplanes is $2 \binom{N}{d}$ and we have $\Pi_{\mathcal{H}}(N) \leq 2^{d+1} \binom{N}{d} < 2^{d+1} \frac{N^d}{d!} = \mathcal{O}(N^d)$.

In addition, there is an algorithm that enumerates all the subsets S_h in $\binom{N}{d}$ iterations and builds \mathcal{H}_S by computing a hyperplane passing through the points¹ in S_h and the corresponding 2^{d+1} labelings at each iteration. Since these inner computations can be performed in constant time with respect to N , the algorithm has a time complexity in the order of $\binom{N}{d} = \mathcal{O}(N^d)$. \square

4.2 Global optimization of binary PWA models

We can use the results above to reduce the complexity of Problem 1 in the binary case, considered in the following in its equivalent form (8) from Proposition 2. First, note that the cost function in (8) only depends on h , since all feasible values of \mathbf{w} for a given h yield the same cost. Furthermore, the cost does not depend on the exact value of h , but only on the resulting classification, i.e., on $h(\mathbf{x}_i)$, $i = 1, \dots, N$. Thus, given a global solution h^* to (8), any classifier h producing the same classification yields the same cost function value and hence is also a global solution. Thus, the problem reduces to the search for the correct classification $\mathbf{q} \in \mathcal{H}_S$, whose complexity is in $\mathcal{O}(\Pi_{\mathcal{H}}(N))$ and bounded by Theorem 3. In addition, for the purpose of binary PWA regression, opposite labelings \mathbf{q} and $-\mathbf{q}$ are equivalent and can be pruned from \mathcal{H}_S . This is due to the symmetry of the cost function (8). Algorithm 1 provides a solution to Problem 1 for the binary case while taking this symmetry into account.

¹The normal \mathbf{h} of a hyperplane $\{\mathbf{x} : \mathbf{h}^T \mathbf{x} + b = 0\}$ passing through d points $\{\mathbf{x}_i\}_{i=1}^d$ in \mathbb{Q}^d can be computed as a unit vector in the null space of $[\mathbf{x}_2 - \mathbf{x}_1, \dots, \mathbf{x}_d - \mathbf{x}_1]^T$, while the offset is given by $b = -\mathbf{h}^T \mathbf{x}_i$ for any of the \mathbf{x}_i 's.

Algorithm 1 Exact solution to Problem 1 for $n = 2$

Input: A data set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N \subset (\mathbb{Q}^d \times \mathbb{Q})^N$.

Initialize $S \leftarrow \{\mathbf{x}_i\}_{i=1}^N$ and $J^* \leftarrow +\infty$.

for all $S_h \subset S$ such that $|S_h| = d$ **do**

 Compute the parameters $(\mathbf{h}_{S_h}, b_{S_h})$ of a hyperplane passing through the points in S_h .

 Classify the data points: $S_1 = \{\mathbf{x}_i \in S : \mathbf{h}_{S_h}^T \mathbf{x}_i + b_{S_h} > 0\}$, $S_2 = \{\mathbf{x}_i \in S : \mathbf{h}_{S_h}^T \mathbf{x}_i + b_{S_h} < 0\}$.

for all classification of S_h into S_h^1 and S_h^2 **do**

 Set $\mathbf{w}_j \in \operatorname{argmin}_{\mathbf{v} \in \mathbb{Q}^{d+1}} \sum_{\mathbf{x}_i \in S_j \cup S_h^j} \ell(y_i - \mathbf{v}^T \bar{\mathbf{x}}_i)$, $j = 1, 2$,

$$J = \frac{1}{N} \sum_{j=1}^2 \sum_{\mathbf{x}_i \in S_j \cup S_h^j} \ell(y_i - \mathbf{w}_j^T \bar{\mathbf{x}}_i),$$

 and update the best solution ($J^* \leftarrow J$, $\mathbf{w}^* \leftarrow [\mathbf{w}_1^T, \mathbf{w}_2^T]^T$, $\mathbf{h}^* \leftarrow \mathbf{h}_{S_h}$, $b^* \leftarrow b_{S_h}$) if $J < J^*$.

end for

end for

return $\mathbf{w}^*, \mathbf{h}^*, b^*$.

Proposition 4. Under Assumptions 1–2, Algorithm 1 exactly solves Problem 1 for $n = 2$ and any fixed d with a polynomial complexity in the order of $T(N)\mathcal{O}(N^d)$.

Proof. By following a similar path as for Theorem 3, Algorithm 1 can be proved to test all linear classifications of the data points up to symmetric ones. Since Algorithm 1 computes a solution in terms of \mathbf{w} that is feasible for (8) for each of these classifications, the value of J coincides with the cost function of (8) for a particular h . By the symmetry of this cost function with respect to h and the fact that it only depends on h via its values at the data points, Algorithm 1 computes all possible values of the cost function, including the exact global optimum of (8), and returns a global minimizer. Thus, by Proposition 2, it also solves Problem 1. The total number of iterations of Algorithm 1 is $2^d \binom{N}{d} = \mathcal{O}(N^d)$ and, under Assumption 2, these iterations only involve operations computed in polynomial time in the order of $T(N)$, hence the overall time complexity in the order of $T(N)\mathcal{O}(N^d)$. \square

4.3 Multi-class extension

For $n > 2$, the boundary between 2 modes j and $k > j$ implemented by a linear classifier from \mathcal{H} in (1) is a hyperplane of equation $h_{jk}(\mathbf{x}) = h_j(\mathbf{x}) - h_k(\mathbf{x}) = 0$, i.e., based on the difference of the two functions $h_j(\mathbf{x}) = \mathbf{h}_j^T \mathbf{x} + b_j$ and $h_k(\mathbf{x}) = \mathbf{h}_k^T \mathbf{x} + b_k$. Based on these hyperplanes, the classification rule can be written as

$$h(\mathbf{x}) = \operatorname{argmax}_{k \in \mathcal{Q}} h_k(\mathbf{x}) = j, \text{ such that } \begin{cases} h_{jk}(\mathbf{x}) \geq 0, \forall k > j, \\ h_{kj}(\mathbf{x}) < 0, \forall k < j. \end{cases}$$

Based on these facts, we can build an algorithm to recover all possible classifications consistent with a linear classification in the sense of (1).

Theorem 4. For the set of multi-class linear classifiers of \mathbb{Q}^d , \mathcal{H} in (1), the growth function is bounded for any $N > d$ by

$$\Pi_{\mathcal{H}}(N) \leq \left[2^{d+1} \binom{N}{d} \right]^{n(n-1)/2} = \mathcal{O}(N^{dn(n-1)/2})$$

and, for any set S of N points of \mathbb{Q}^d in general position, an algorithm builds \mathcal{H}_S in $\mathcal{O}(N^{dn(n-1)/2})$ time.

Proof. Any classification produced by a classifier from (1) can be computed from the signs of the $n_H = n(n-1)/2$ functions $h_{jk} = h_j - h_k$, $1 \leq j < k \leq n$, corresponding to the pairwise separating hyperplanes. For any S , for each of these hyperplanes, Proposition 3 provides an equivalent binary classifier which must be one from the $2 \binom{N}{d}$ hyperplanes passing through d points S_{jk} of S . The number of sets of n_H such hyperplanes is $2^{n_H} \binom{N}{d}^{n_H}$. Since these classifiers cannot produce all the $2^{n_H d}$ classifications of the $n_H d$ points in the sets S_{jk} , we must also take these into account so that the number of classifications of S is upper bounded by $|\mathcal{H}_S| \leq 2^{n_H d} 2^{n_H} \binom{N}{d}^{n_H} = \mathcal{O}(N^{dn_H})$. This upper bound holds for any S , and thus also applies to the growth function. Finally, an algorithm that makes explicit all the classifications mentioned above to build \mathcal{H}_S can be constructed in a recursive manner, with one classification per iteration and thus with a similar number of iterations, each one including computations performed in constant time. \square

Theorem 4 implies the following for PWA regression.

Corollary 1. *Under Assumptions 1–2, a global solution to Problem 1 with $n \geq 3$ can be computed with a polynomial complexity in the order of $T(N)\mathcal{O}(N^{dn(n-1)/2})$.*

5 Conclusions

The paper discussed complexity issues for PWA regression and showed that i) the global minimization of the error is NP-hard in general, and ii) for fixed number of modes and data dimension, an exact solution can be obtained in time polynomial in the number of data. The proof of NP-hardness also implies that the problem remains NP-hard even when the number of modes is fixed to 2, which indicates that the complexity is mostly due to the data dimension. An open issue concerns the conditions under which a PWA system generates trajectories satisfying the general position assumption used by the polynomial-time algorithm. Future work will also focus on the extension of the results to the case of arbitrarily switched systems and heuristics inspired by the polynomial-time algorithm, whose practical application remains limited by an exponential complexity in the dimension.

Acknowledgements

The author would like to thank the anonymous reviewers for their comments and suggestions. Thanks are also due to Yann Guermeur for carefully reading this manuscript.

References

- [1] S. Paoletti, A. L. Juloski, G. Ferrari-Trecate, R. Vidal, Identification of hybrid systems: a tutorial, *European Journal of Control* 13 (2-3) (2007) 242–262.
- [2] A. Garulli, S. Paoletti, A. Vicino, A survey on switched and piecewise affine system identification, in: *Proc. of the 16th IFAC Symp. on System Identification (SYSID)*, 2012, pp. 344–355.
- [3] R. Vidal, S. Soatto, Y. Ma, S. Sastry, An algebraic geometric approach to the identification of a class of linear hybrid systems, in: *Proc. of the 42nd IEEE Conf. on Decision and Control (CDC)*, Maui, Hawaiï, USA, 2003, pp. 167–172.
- [4] L. Bako, Identification of switched linear systems via sparse optimization, *Automatica* 47 (4) (2011) 668–677.
- [5] I. Maruta, H. Ohlsson, Compression based identification of PWA systems, in: *Proc. of the 19th IFAC World Congress*, Cape Town, South Africa, 2014, pp. 4985–4992.

- [6] L. Ljung, System identification: Theory for the User, 2nd Edition, Prentice Hall, 1999.
- [7] G. Ferrari-Trecate, M. Muselli, D. Liberati, M. Morari, A clustering technique for the identification of piecewise affine systems, *Automatica* 39 (2) (2003) 205–217.
- [8] A. Bemporad, A. Garulli, S. Paoletti, A. Vicino, A bounded-error approach to piecewise affine system identification, *IEEE Transactions on Automatic Control* 50 (10) (2005) 1567–1580.
- [9] A. L. Juloski, S. Weiland, W. Heemels, A Bayesian approach to identification of hybrid systems, *IEEE Transactions on Automatic Control* 50 (10) (2005) 1520–1533.
- [10] F. Lauer, G. Bloch, R. Vidal, A continuous optimization framework for hybrid system identification, *Automatica* 47 (3) (2011) 608–613.
- [11] F. Lauer, Estimating the probability of success of a simple algorithm for switched linear regression, *Nonlinear Analysis: Hybrid Systems* 8 (2013) 31–47, supplementary material available at <http://www.loria.fr/~lauer/klinreg/>.
- [12] T. Pham Dinh, H. Le Thi, H. Le, F. Lauer, A difference of convex functions algorithm for switched linear regression, *IEEE Transactions on Automatic Control* 59 (8) (2014) 2277–2282.
- [13] N. Ozay, M. Sznaiar, C. Lagoa, O. Camps, A sparsification approach to set membership identification of switched affine systems, *IEEE Transactions on Automatic Control* 57 (3) (2012) 634–648.
- [14] H. Ohlsson, L. Ljung, S. Boyd, Segmentation of ARX-models using sum-of-norms regularization, *Automatica* 46 (6) (2010) 1107–1111.
- [15] H. Ohlsson, L. Ljung, Identification of switched linear regression models using sum-of-norms regularization, *Automatica* 49 (4) (2013) 1045–1050.
- [16] F. Lauer, V. L. Le, G. Bloch, Learning smooth models of nonsmooth functions via convex optimization, in: *Proc. of the IEEE Int. Workshop on Machine Learning for Signal Processing (MLSP)*, Santander, Spain, 2012.
- [17] J. Roll, A. Bemporad, L. Ljung, Identification of piecewise affine systems via mixed-integer programming, *Automatica* 40 (1) (2004) 37–50.
- [18] M. Garey, D. Johnson, *Computers and Intractability: a Guide to the Theory of NP-Completeness*, W.H. Freeman and Company, 1979.
- [19] D. Aloise, A. Deshpande, P. Hansen, P. Popat, NP-hardness of Euclidean sum-of-squares clustering, *Machine Learning* 75 (2) (2009) 245–248.
- [20] V. Blondel, J. Tsitsiklis, A survey of computational complexity results in systems and control, *Automatica* 36 (9) (2000) 1249–1274.
- [21] R. Alur, N. Singhanian, Precise piecewise affine models from input-output data, in: *Proc. of the 14th Int. Conf. on Embedded Software (EMSOFT)*, 2014.
- [22] G. H. Golub, C. F. Van Loan, *Matrix Computations*, 4th Edition, John Hopkins University Press, 2013.
- [23] V. Vapnik, *Statistical Learning Theory*, John Wiley & Sons, 1998.